

(19) World Intellectual Property Organization
International Bureau



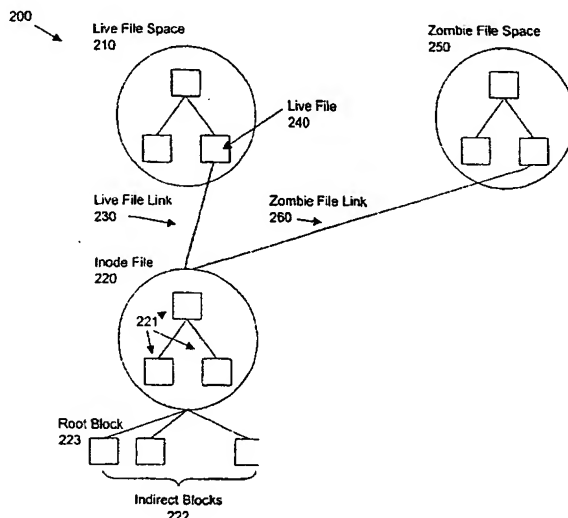
(43) International Publication Date
7 March 2002 (07.03.2002)

PCT

(10) International Publication Number
WO 02/19110 A2

- (51) International Patent Classification⁷: G06F 11/00 (74) Agent: SWERNOFSKY, Steven, A.; Swernofsky Law Group, P.O. Box 390013, Mountain View, CA 94039-0013 (US).
- (21) International Application Number: PCT/US01/25901
- (22) International Filing Date: 17 August 2001 (17.08.2001) (84) Designated States (regional): European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR).
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data: 09/642,066 18 August 2000 (18.08.2000) US
- (71) Applicant: NETWORK APPLIANCE, INC. [US/US]; 495 East Java Drive, Sunnyvale, CA 94089 (US).
- (72) Inventors: CHEN, Ray; 400 Castro Court, Campbell, CA 95008 (US). EDWARDS, John, K.; 1173 Crandano Court, Sunnyvale, CA 94087-2076 (US). PATEL, Kayuri; 20380 Stevens Creek Blvd., Apt. 219, Cupertino, CA 95014 (US).
- Declaration under Rule 4.17:**
— as to the applicant's entitlement to claim the priority of the earlier application (Rule 4.17(iii)) for all designations
- Published:**
— without international search report and to be republished upon receipt of that report
- For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.*

(54) Title: MANIPULATION OF ZOMBIE FILES AND EVIL-TWIN FILES



(57) Abstract: The invention provides a method and system for reliably performing extra-long operations in a reliable state-full system (such as a file system). The system records consistency points, or otherwise assures reliability, notwithstanding the continuous performance of extra-long operations and the existence of intermediate states for those extra-long operations. Moreover, performance of extra-long operations is both deterministic and atomic with regard to consistency points (or other reliability techniques used by the system). The file system includes a separate portion of the file system reserved for files having extra-long operations in progress, including file deletion and file truncation. This separate portion of the file system is called the zombie filesystem; it includes a separate name space from the regular ("live") file system that is accessible to users, and is maintained as part of the file system when recording a consistency point. The file system includes a file deletion manager that determines, before beginning any file deletion operation, whether it is necessary to first move the file being deleted to the zombie filesystem. The file system includes a zombie file deletion manager that performs portions of

the file deletion operation on zombie files in atomic units. The file system also includes a file truncation manager that determines, before beginning any file truncation operation, whether it is necessary to create a complementary file called an "evil twin". The truncation manager will move all blocks to be truncated from the file being truncated to the evil twin file. The file system includes a zombie file truncation manager that performs portions of the file truncation operation on the evil-twin file in atomic units. An additional advantage provided by the file system is that files having attached data elements, called "composite" files, can be subject to file deletion and other extra-long operations in a natural and reliable manner. The file system moves the entire composite file to the zombie filesystem, deletes each attached data element individually, and thus resolves the composite file into a non-composite file. If the non-composite file is sufficiently small, the file deletion manager can delete the non-composite file without further need for the zombie filesystem. However, if the non-composite file is sufficiently large, the file deletion manager can delete the non-composite file using the zombie filesystem.

WO 02/19110 A2

MANIPULATION OF ZOMBIE FILES AND EVIL-TWIN FILES

Background of the Invention5 1. *Field of the Invention*

This invention relates to file server systems, including those file server systems in which it is desired to maintain reliable file system consistency.

10 2. *Related Art*

In systems providing file services, such as those including file servers and similar devices, it is generally desirable for the server to provide a file system that is reliable despite the possibility of error. For example, it is desirable to provide a file system that is
15 reliably in a consistent state, regardless of problems that might have occurred with the file server, and regardless of the nature of the file system operations requested by client devices.

One known method of providing reliability in systems that maintain state (including such state as the state of a file system or other set of data structures) is to provide
20 for recording checkpoints at which the system is known to be in a consistent state. Such checkpoints, sometimes called "consistency points," each provide a state to which the system can retreat in the event that an error occurs. From the most recent consistency point, the system can reattempt each operation to reach a state it was in before the error.

25 One problem with this known method is that some operations can require substantial amounts of time in comparison with the time between consistency points. For example, in the WAFL file system, operations on very large files can require copying or modifying very large numbers of file blocks in memory or on disk, and can therefore take a substantial fraction of the time from one consistency point to another. In the WAFL file
30 system, two such operations are deleting very large files and truncating very large files. Accordingly, it might occur that recording a consistency point cannot occur properly while one of these extra-long operations is in progress.

The fundamental requirement of a reliable file system is that the state of the file system recorded on non-volatile storage must reflect only completed file system operations. In the case of a file system like WAFL that issues checkpoints, every file system operation must be complete between two checkpoints. In the earliest versions of the WAFL file system there was no file deletion manager present, thus very large files created a problem as it was possible that such large files could not be deleted between the execution of two consistency checkpoints.

This problem was partially solved in later versions of the WAFL file system, where a file deletion manager was assigned to perform the operation of file deletion, and a consistency point manager was assigned to perform the operation of recording a consistency point. The file deletion manager would attempt to resolve the problem of extra-long file deletions by repeatedly requesting more time from the consistency point manager, thus "putting off" the consistency point manager until a last-possible moment. However, at that last-possible moment, the file deletion manager would be required to give way to the consistency point manager, and allow the consistency point manager to record the consistency point. When this occurred, the file deletion manager would be unable to complete the file deletion operation. In that earlier version of the WAFL file system, instead of completing the file deletion operation, the file deletion manager would move the file to a fixed-length "zombie file" list to complete the file deletion operation. At a later time, a zombie file manager would re-attempt the file deletion operation for those files on the fixed-length zombie file list.

While this earlier method achieved the general result of performing file deletions on very large files, it has the drawbacks that it is a source of unreliability in the file system. First, the number of files that could be processed simultaneously as zombie files was fixed in the previous version.

Second, the file deletion manager and crash recovery mechanism did not communicate. The file deletion manager did not notify the crash recovery mechanism that a file was being turned into a zombie and the crash recovery mechanism was unable to create zombie files. Thus, to allow a checkpoint to be recorded, a long file would have to be turned into a zombie. If the system crashed at this point, the crash recovery mechanism might not be able to correctly recover the file system since it is unaware that a zombie file should be

created and was incapable of creating zombie files should the need arise. Similarly, the operations of the file deletion manager when creating zombie files, and its operations in deleting those zombie files, were not recorded in non-volatile storage, and thus could not be "replayed" after recovery to duplicate the operations of the file deletion manager.

5

Third, since the file deletion manager and replay mechanism did not communicate, the free space reported could be inaccurately reported. Attempts to restore state could fail, because the amount of free space could be different than that actually available. Attempts to restore state could also fail because the operations of the file deletion manager in using zombie files were not recorded in non-volatile storage; as a result, it might occur that other operations performed during replay could conflict with the file deletion manager and cause a crash.

Fourth, the earlier method is non-deterministic in the sense that it is not assured whether any particular file deletion operation will be completed before or after a selected consistency point. Moreover, the earlier method does not resolve problems associated with other extra-long file operations, such as requests to truncate very large files to much smaller length.

Accordingly, it would be advantageous to provide a technique for extra-long operations in a reliable state-full system (such as a file system) that is not subject to the drawbacks of the known art. Preferably, in such a technique, those parts of the system responsible for recording of consistency points are fully aware of the intermediate states of extra-long operations, the performance of extra-long operations is relatively deterministic, and performance of extra-long operations is atomic with regard to consistency points.

Summary of the Invention

The invention provides a method and system for reliably performing extra-long operations in a reliable state-full system (such as a file system). The system records consistency points, or otherwise assures reliability (such as using a persistent-memory log file), notwithstanding the continuous performance of extra-long operations and the existence of intermediate states for those extra-long operations. The system provides for replay, after recovery, of those portions of extra-long operations which were completed, thus assuring

that recovery and replay are consistent with operations of the file deletion manager and the zombie deletion manager. Moreover, performance of extra-long operations is both deterministic and atomic with regard to consistency points (or other reliability techniques used by the system).

5

The file system includes a separate portion reserved for files having extra-long operations in progress, including file deletion and file truncation;. this separate portion of the file system is called the zombie filesystem. The zombie filesystem includes a separate name space from the regular ("live") file system and is maintained as part of the file system when recording a consistency point, just like the live filesystem. The live filesystem refers to those files that are accessible to users in normal operation, such as for example those files for which a path can be traced from a root of a hierarchical namespace. The file system includes a file deletion manager that determines, before beginning any file deletion operation, whether it is necessary to first move the file being deleted to the zombie filesystem. The file system includes a zombie file deletion manager that performs portions of the file deletion operation on zombie files in atomic units.

The file system also includes a file truncation manager. Before beginning any file truncation operation, the file truncation manager determines whether it is necessary to create a complementary file called an "evil twin" file, located in the zombie filesystem. The truncation manager will move all blocks to be truncated from the file being truncated to the evil twin file. Moving blocks is typically faster and less resource-intensive than deleting blocks. The "evil twin" is subsequently transformed into a zombie file. The file system includes a zombie file truncation manager that can then perform truncation of the zombie file asynchronously in atomic units. Furthermore, the number of files that can be linked to the zombie filesystem is dynamic, allowing the zombie filesystem the ability to grow and shrink as required to process varying numbers of files.

An additional advantage provided by the file system is that files having attached data elements, called "composite" files, can be subject to file deletion and other extra-long operations in a natural and reliable manner. When performing such operations for composite files, the file system moves the entire composite file to the zombie filesystem, deletes each attached data element individually, and thus resolves the composite file into a non-composite file. If the non-composite file is sufficiently small, the file deletion manager

can delete the non-composite file without further need for the zombie filesystem. However, if the non-composite file is sufficiently large, the file deletion manager can delete the non-composite file using the zombie filesystem.

5 The invention provides an enabling technology for a wide variety of applications for reliable systems, so as to obtain substantial advantages and capabilities that are novel and non-obvious in view of the known art. Examples described below primarily relate to reliable file systems, but the invention is broadly applicable to many different types of systems in which reliability and extra-long operations are both present.

10

Brief Description of the Drawings

Figure 1 shows a block diagram of a portion of a system using a zombie filesystem.

15

Figure 2 illustrates a file structure in a system using a zombie filesystem.

Figure 3 shows a process flow diagram for file deletion in a method for operating a system for manipulation of zombie files and evil-twin files.

20

Figure 4 shows a process flow diagram for file truncation in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

25 Figure 5 shows a process flow diagram for replaying operations in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

Lexicography

30 The following terms refer to or relate to aspects of the invention as described below. The descriptions of general meanings of these terms are not intended to be limiting, only illustrative.

- live filesystem — This term generally refers to a portion of the file system where files are available to users in normal operation. In a preferred embodiment, the live filesystem includes those inodes (or other types of file control structure) that are not yet allocated to in-use files.

5

- zombie filesystem — This term generally refers to a portion of the file system where files are not available to users in normal operation, but can still be manipulated by the file system as if they were normal files.

10

- Storage Operating System — in general refers to the computer-executable code operable on a storage system that implements file system semantics and manages data access. In this sense, ONTAP software is an example of such a storage operating system implemented as a microkernel, with its WAFL layer implementing the file system semantics. The storage operating system can also be implemented as an application program operating over a general-purpose operating system, such as UNIX® or Windows NT®, or as a general-purpose operating system with configurable functionality, which is configured for storage applications.

15

As noted above, these descriptions of general meanings of these terms are not intended to be limiting, only illustrative. Other and further applications of the invention, including extensions of these terms and concepts, would be clear to those of ordinary skill in the art after perusing this application. These other and further applications are part of the scope and spirit of the invention, and would be clear to those of ordinary skill in the art, without further invention or undue experimentation.

25

Detailed Description of the Preferred Embodiment

In the following description, a preferred embodiment of the invention is described with regard to preferred process steps and data structures. Embodiments of the invention can be implemented using general-purpose processors or special purpose processors operating under program control, or other circuits, adapted to particular process steps and data structures described herein. Implementation of the process steps and data structures described herein would not require undue experimentation or further invention.

30

System Elements

Figure 1 shows a block diagram of a portion of a system using a zombie
5 filespace.

A system 100 includes a file server 110 including a processor 111, program
and data memory 112, a network interface card 115, and mass storage 120.

10 The program and data memory 112 include program instructions and data
structures used by a file deletion manager 121, a zombie file deletion manager 122, a file
truncation manager 123, or a zombie file truncation manager 124.

The file deletion manager 121 responds to a file server request (such as one
15 received from a user of the file server 110), and performs an operation for deleting a file. As
shown herein, the operation for deleting a file might include transferring the file from a live
filespace 210 (shown in figure 2) to a zombie filespace 250 (shown in figure 2) and
performing additional operations on the file in the zombie filespace 250. The zombie file
deletion manager 122 performs these additional operations.

20 Similarly, the file truncation manager 123 responds to a file server request
(such as one received from a user of the file server 110), and performs an operation for
deleting a file. As shown herein, the operation for deleting a file might include transferring
the file to a zombie filespace 250 and performing additional operations on the file in the
25 zombie filespace 250. The zombie file truncation manager 124 performs these additional
operations.

The network interface card 115 couples the file server 110 to a network. In a
preferred embodiment, the network includes an Internet, intranet, extranet, virtual private
30 network, enterprise network, or another form of communication network.

The mass storage 120 can include any device for storing relatively large
amounts of information, such as magnetic disks or tapes, optical drives, or other types of
mass storage.

File Structure Example

Figure 2 illustrates a file structure in a system using a zombie filesystem.

5

A file structure 200 includes, a live filesystem 210, an inode file 220, a live file link 230, a file 240, a zombie filesystem 250, and a zombie file link 260.

10 The live filesystem 210 contains a live root block 211 and all associated blocks of data and metadata for live files. As noted above, "live files" are files in the live filesystem, which may be accessed by users in normal operation.

The inode file 220 is associated with the file to be deleted and contains information about the file. The inode file 220 itself is preferably recorded using a tree structure, in which individual entries 221 for files (including their live file links 230) are maintained at leaves of the tree, and in which one or more indirect blocks 222 are maintained at nodes of the tree to allow the entire inode file 220 to be reached from a root block 223 therefor. Small inode files 220 might not require any indirect blocks 222, or might even be stored directly in data blocks for their containing directory.

20

The live file link 230, links a file to the live filesystem 210.

Similar to an inode file 220, the file 240 includes a plurality of file blocks 241, and a plurality of block links 242. The file blocks 241 are connected by the plurality of block links 242. The file 240 is illustrative of a file to be deleted. The structure of the file as defined above is a hierarchical tree-like structure, however, there is no requirement in any embodiment of the invention that the invention be applied only to file structures (or inode structures) of this type. The use of a hierarchical tree-like structure filing system is intended to be illustrative only and not limiting.

30

If the file is a composite file, it has attached data elements 243 which are associated with the file 240 (such as possibly by one or more references from the file's inode file 220).

The zombie filesystem 250 contains a zombie root block 251 and all associated blocks of data for zombie files (files in the zombie filesystem, which are in the process of being deleted or truncated).

5 The zombie file link 260 links a file to be deleted to the zombie filesystem 250. A file that has been linked to the zombie filesystem 250 is referred to as a "zombie file" while it is so linked. Zombie files in the zombie filesystem 250 are maintained in like manner as live files 240 in the live filesystem 210.

10 *Method of Operation — File Deletion*

Figure 3 shows a process flow diagram for file deletion in a method for operating a system for manipulation of zombie files and evil-twin files.

15 A method 300 includes a set of flow points and a set of steps. The system 100 performs the method 300. Although the method 300 is described serially, the steps of the method 300 can be performed by separate elements in conjunction or in parallel, whether asynchronously, in a pipelined manner, or otherwise. There is no particular requirement that the method 300 be performed in the same order in which this description lists the steps,
20 except where so indicated.

In this method 300, each operation denoted by a flow point is recorded in a file system log, such as a persistent memory that can be accessed in the event of a file system crash or other service interruption. The file system can and does generate checkpoints while
25 these operations are being performed. After a crash, the file system replays the operations noted in the log, as further described with regard to figure 5.

At a flow point 310, a system user selects the file 240 for deletion. User interfaces for this activity vary from system to system but are well known in the art.

30

At a flow point 320, the file 240 is identified by the system as a large file requiring zombie processing. In a preferred embodiment, the specific size of a file necessary to trigger zombie processing is parameter-based, software-selectable, however, it can be any

set of instructions supporting this functionality, such as instructions hard-coded on a computer chip.

5 The file 220 is identified as a large file in response to an amount of time calculated as necessary to delete the file 220. The amount of time is calculated in response to a number of data blocks included in the file, and in response to a size on record for the file. In a preferred embodiment, the file 220 is identified as a large file if it has more than one indirect block 241, that is, if the file 220 has more than about 1,024 data blocks 241. In a preferred embodiment, all composite files 220 are also identified as large files for this
10 purpose.

In alternative embodiments, depending on the underlying implementation of the file system and storage operating system, the file is identified as a large file in response to other metrics of when extra-long operations can consume too many resources at once,
15 hold resources locked for too long a period of time, or otherwise consume too much of a single resource, or some combination thereof, so as to jeopardize correct operation of other parts of the file system and storage operating system. Examples of such other metrics include an amount of log space, a number of log entries, or some other measure of unfinished work needed to be completed, that would be used if the deletion (or truncation)
20 operation is too large.

At a flow point 325, the file deletion manager 121 determines whether the zombie filesystem 250 needs to be enlarged to accommodate another zombie file, and if necessary enlarges the zombie filesystem.

25

In a preferred embodiment, the file deletion manager 121 attempts to allocate an entry in the zombie filesystem 250. If this is possible (that is, at least one entry is available in the zombie filesystem 250 for use), the file deletion manager 121 can proceed without requesting enlargement of the zombie filesystem 250. If there is no entry available in the
30 zombie filesystem 250 for use, the file deletion manager 121 requests the file server 110 to enlarge the zombie filesystem 250 (such as by creating another free entry therein), and proceeds to allocate the newly created free entry for use. If the newly created free entry has

been allocated by another process, the file deletion manager 121 repeats this flow point until it is able to allocate an entry for its own use.

At a flow point 330, the link connecting the file 240 to the live filespace 210 is terminated. At this point the file 240 is no longer available to users connected to the file server 110.

In a preferred embodiment, the file deletion manager 121 also alters the generation number of the inode 220 for the file 210, so that external users of the file server 110 can no longer refer to the file 210 by file handles they might have kept. Those users will see the file 210 as having disappeared (that is, been deleted).

At a flow point 340, the file 240 is linked to the zombie filespace 250 via the zombie file link 260. At this point, file 240 is referred to as a zombie file.

At a flow point 350, the zombie file deletion manager 122 starts deleting portions of the file 240 by terminating block links 242 at the outer leaves of the file tree. As file blocks 241 are deleted by the zombie deletion manager 122, they become available for storage of other data. This fact is reflected in the free space indicator of the mass storage 120.

At a flow point 360, the file 240 is deleted. Since the file 240 in the zombie filespace 250 has been deleted, this is equivalent to freeing the inode 220, and any other file system control structure, for the file 240, and terminating any link between the file 240 and the zombie filespace 250.

Method of Operation - File Truncation

Figure 4 shows a process flow diagram for file truncation in a method for operating a system Manipulation of Zombie Files and Evil-Twin Files.

A method 400 includes a set of flow points and a set of steps. The system 100 performs the method 400. Although the method 400 is described serially, the steps of the method 400 can be performed by separate elements in conjunction or in parallel, whether

asynchronously, in a pipelined manner, or otherwise. There is no particular requirement that the method 400 be performed in the same order in which this description lists the steps, except where so indicated.

5 In this method 400, each operation denoted by a flow point is recorded in a file system log, such as a persistent memory that can be accessed in the event of a file system crash or other service interruption. The file system can and does generate checkpoints while these operations are being performed. After a crash, the file system replays the operations noted in the log, as further described with regard to figure 5.

10

At a flow point 410, a system user selects the file 240 for truncation. User interfaces for this activity vary from system to system but are well known in the art.

15 At a flow point 420, the file system (that is, the file system component of the storage operating system) identifies the amount of the file to be truncated as requiring evil twin/zombie processing. In the preferred embodiment, the specific amount of data to be truncated necessary to trigger evil twin/zombie processing is parameter-based software-selectable; however, it can be any set of instructions supporting this functionality, such as instructions hard-coded on a computer chip. In a preferred embodiment, identification of a
20 file for evil twin processing is similar to identification of a file for zombie processing.

At a flow point 425, the file truncation manager 123 determines whether the zombie filesystem 250 needs to be enlarged to accommodate another zombie file, and if necessary enlarges the zombie filesystem. This flow point is similar to the flow point 325.

25

At a flow point 430, an evil twin file is created. At this point the file 240 is unavailable to the user. This flow point is similar to the flow points 330 and 340, except that the original file is not removed from the live filesystem 210.

30

At a flow point 440, blocks of data to be truncated are moved from the file 240 to the evil twin file. Links associating the data blocks to be truncated from the live file in the live filesystem are broken, and corresponding links associating the same data blocks with the evil twin file in the zombie filesystem are created. This flow point is similar to the

flow points 330 and 340, except that only a subset of the data blocks in the original file are removed from the live filesystem 210 and transferred to the zombie filesystem 250.

At a flow point 450, file attributes for the file 240 are adjusted appropriately
5 (for example, the size of the file, the number of blocks in the file, and the file's timestamp).

At a flow point 460, the evil twin file is turned into a zombie file. It is connected to the zombie filesystem. This flow point is similar to the flow point 340, except that it is the evil twin, not the original file, which is linked to the zombie filesystem 250.

10

At a flow point 470, the file 240 is marked as available in the live filesystem. At this point the file 240 is available to all users.

At a flow point 480, the zombie file deletion manager 122 frees all blocks
15 attached to the zombie file.

At a flow point 490, the zombie file has been deleted and the link to the zombie filesystem 250 is terminated. Since the zombie file in the zombie filesystem 250 has been deleted, this is equivalent to freeing the inode 220, and any other file system control
20 structure, for the zombie file.

Method of Operation — Replay

Figure 5 shows a process flow diagram for replaying operations in a method
25 for operating a system Manipulation of Zombie Files and Evil-Twin Files.

A method 500 includes a set of flow points and a set of steps. The system 100 performs the method 500. Although the method 500 is described serially, the steps of the method 500 can be performed by separate elements in conjunction or in parallel, whether
30 asynchronously, in a pipelined manner, or otherwise. There is no particular requirement that the method 500 be performed in the same order in which this description lists the steps, except where so indicated.

At a flow point 510, the file server 110 has recovered from a crash or other service interruption.

At a step 511, the file server 110 examines its log (preferably recorded in a persistent memory), and determines which log entries should be replayed. In a preferred embodiment, those log entries not marked in the log as being committed as part of a consistency point are required to be replayed. In a preferred embodiment, the log is recorded in a persistent memory and pointed to by at least one link from a persistently recorded file system control block. To quickly determine this, the file system control block is preferably
5
10
15
flagged as being "clean" when the system is shut down normally. When rebooting, the system can check each file system to determine if was shut down cleanly. If it was not, then log entries that reflect changes not present in the on-disk file system must be replayed. There are known techniques for determining which such log entries. One method is time-stamping when log entries and the file system control block were last updated.

At a step 512, the file server 110 replays the operation designated by each log entry, thus re-performing those operations.

At an optional (but preferred) step 513, the file server 110 generates a
20
checkpoint when all log entries have been replayed.

At a flow point 520, the file server 110 has both recovered from the crash or other service interruption, and replayed all necessary log entries, so normal file server operations can proceed.

Generality of the Invention

The invention has general applicability to various fields of use, not necessarily related to the services described above. For example, these fields of use can
25
30
include one or more of, or some combination of, the following:

- The invention is applicable to all computer systems utilizing large files.

- The invention is applicable to all computer systems performing long-duration operations on files.

Other and further applications of the invention in its most general form, will
5 be clear to those skilled in the art after perusal of this application, and are within the scope
and spirit of the invention.

Although preferred embodiments are disclosed herein, many variations are
possible which remain within the concept, scope, and spirit of the invention, and these
10 variations would become clear to those skilled in the art after perusal of this application.

Claims

1. A method of operating a filesystem, said filesystem including a live
filesystem accessible to users and a zombie filesystem not accessible to users, said method
5 including

recording changes to said zombie filesystem in a persistent memory.

2. A method as in claim 1, including, for a deletion operation on a file in
said live filesystem,

10 transferring said file from said live filesystem to said zombie filesystem;
breaking links associating disk blocks with said file in a plurality of steps
while said file is associated with said zombie filesystem, wherein said recording of changes
includes recording said breaking of links in a plurality of steps; and
altering said live filesystem to reflect said deletion operation.

15 3. A method as in claim 1, including, for a truncation operation on a file
in said live filesystem,

transferring at least a portion of said file from said live filesystem to said
zombie filesystem;
20 breaking links associating disk blocks with said file in a plurality of steps
while a portion of said file is associated with said zombie filesystem, wherein said recording
of changes includes recording said breaking of links in a plurality of steps; and
altering said live filesystem to reflect changes associated with said breaking of
links.

25 4. A method as in claim 1, including, for an operation apparent to users
as substantially atomic, performing said operation in a plurality of steps using said zombie
filesystem, wherein said recording changes is performed in said persistent memory for each of
said plurality of steps.

30 5. A method as in claim 1, including, for an operation performed on a
file having attached data elements, performing said operation using said zombie filesystem.

6. A method as in claim 1, including, for an operation performed using said zombie filesystem, altering a size of said zombie filesystem during performance of said operation.

5 7. A method as in claim 1, including, for an operation performed using said zombie filesystem, checkpointing said filesystem during performance of said operation.

8. A method as in claim 1, including recording changes to said live filesystem in said persistent memory, wherein records of changes to said live filesystem and of
10 changes to said zombie filesystem are substantially interspersed.

9. A method as in claim 1, including replaying a set of said changes in response to said record.

15 10. A method as in claim 1, including replaying a set of said changes to said live filesystem and to said zombie filesystem, wherein replay of changes includes substantial interspersed performance of changes to said live filesystem and to said zombie filesystem.

20 11. A method as in claim 1, including replaying a set of said changes in said record in response to a crash recovery by said filesystem.

12. A method as in claim 1, wherein said persistent memory includes a log of substantially all changes, within a selected time duration, to either said live filesystem
25 or said zombie filesystem.

13. A method as in claim 1, wherein said persistent memory includes a log of substantially all changes, within a selected time duration, to said zombie filesystem.

30 14. A method as in claim 1, wherein said recorded changes include a set of substantially atomic operations to said zombie filesystem.

15. A method of operating a filesystem, said filesystem including a live filesystem accessible to users and a zombie filesystem not accessible to users, said method including

dynamically growing said zombie filesystem.

5

16. A method as in claim 15, including, for a deletion or truncation operation on a file in said live filesystem,

allocating storage within said zombie filesystem for metadata associated with said file;

10 performing said dynamic growth in response to failure of said allocation of storage;

re-performing said allocation of storage after said dynamic growth; and transferring said file from said live filesystem to said zombie filesystem.

15 17. A method as in claim 15, wherein said dynamic growth occurs, for an operation performed using said zombie filesystem, during performance of said operation.

18. A method of operating a filesystem, said filesystem including a live filesystem accessible to users and a zombie filesystem not accessible to users, said method including

20

transfer of a file to said zombie filesystem before breakage of links to blocks in said file, in response to an operation on said file, said operation using said zombie filesystem.

19. A method as in claim 18, wherein, for a deletion operation on a file in said live filesystem,

25

said transfer includes

creating a link associating said file with said zombie filesystem; and

breaking a link associating said file with said live filesystem;

and said deletion operation includes

30

breaking links associating disk blocks with said file in a plurality of steps while said file is associated with said zombie filesystem, wherein said recording of changes includes recording said breaking of links in a plurality of steps; and

altering said live filesystem to reflect said deletion operation.

20. A method as in claim 18, wherein, for a truncation operation on a file in said live filesystem,

said transfer includes

creating a link associating at least a portion of said file with said zombie

5 filesystem; and

breaking a link associating said portion with said file in said live filesystem;

and said truncation operation includes

breaking links associating disk blocks with said file in a plurality of steps while a portion of said file is associated with said zombie filesystem, wherein said recording

10 of changes includes recording said breaking of links in a plurality of steps; and

altering said live filesystem to reflect changes associated with said breaking of links.

21. A method of operating a filesystem, said filesystem including a live
15 filesystem accessible to users and a zombie filesystem not accessible to users, said method including

transfer of a file to said zombie filesystem before performing any substantial portion of an operation on said file, said operation using said zombie filesystem.

22. A method as in claim 21, wherein, for a deletion operation on a file in
20 said live filesystem,

said transfer includes

creating a link associating said file with said zombie filesystem; and

breaking a link associating said file with said live filesystem;

25 and said deletion operation includes

breaking links associating disk blocks with said file in a plurality of steps only while said file is associated with said zombie filesystem, wherein said recording of changes includes recording said breaking of links in a plurality of steps; and

altering said live filesystem to reflect said deletion operation.

30 23. A method as in claim 21, wherein, for a truncation operation on a file in said live filesystem,

said transfer includes

creating a link associating at least a portion of said file with said zombie
filesystem; and

breaking a link associating said portion with said file in said live filesystem;
and said truncation operation includes

5 breaking links associating disk blocks with said file in a plurality of steps
only while a portion of said file is associated with said zombie filesystem, wherein said
recording of changes includes recording said breaking of links in a plurality of steps; and
altering said live filesystem to reflect changes associated with said breaking of
links.

10

24. A method of operating a filesystem, said filesystem including a live
filesystem accessible to users and a zombie filesystem not accessible to users, said method
including

replay of an operation on a file, said operation using said zombie filesystem.

15

25. A method as in claim 24, wherein said replay is responsive to a set of
recorded changes in a persistent memory;

and including, for a deletion operation on a file in said live filesystem,

transferring said file from said live filesystem to said zombie filesystem, and

20 recording said transfer in said persistent memory;

breaking links associating disk blocks with said file in a plurality of steps
while said file is associated with said zombie filesystem, and recording said breaking of links
in said persistent memory in a plurality of steps; and

altering said live filesystem to reflect said deletion operation, and recording

25 said alteration in said persistent memory.

26. A method as in claim 24, wherein said replay is responsive to a set of
recorded changes in a persistent memory;

and including, for a truncation operation on a file in said live filesystem,

30 transferring at least a portion of said file from said live filesystem to said
zombie filesystem, and recording said transfer in said persistent memory;

breaking links associating disk blocks with said file in a plurality of steps
while a portion of said file is associated with said zombie filesystem, and recording said
breaking of links in said persistent memory in a plurality of steps; and

altering said live filesystem to reflect changes associated with said breaking of links, and recording said alteration in said persistent memory.

27. A method of operating a filesystem, said filesystem including a live
5 filesystem accessible to users and a zombie filesystem not accessible to users, said method including

replay of a set of filesystem operations, said operations including at least some operations using said live filesystem and at least some operations using said zombie filesystem.

10

28. A method as in claim 27, wherein said replay is responsive to a set of recorded changes in a persistent memory;

and including, for a deletion operation on a file in said live filesystem,

transferring said file from said live filesystem to said zombie filesystem, and
15 recording said transfer in said persistent memory;

breaking links associating disk blocks with said file in a plurality of steps while said file is associated with said zombie filesystem, and recording said breaking of links in said persistent memory in a plurality of steps; and

altering said live filesystem to reflect said deletion operation, and recording
20 said alteration in said persistent memory.

29. A method as in claim 27, wherein said replay is responsive to a set of recorded changes in a persistent memory;

and including, for a truncation operation on a file in said live filesystem,

transferring at least a portion of said file from said live filesystem to said
25 zombie filesystem, and recording said transfer in said persistent memory;

breaking links associating disk blocks with said file in a plurality of steps while a portion of said file is associated with said zombie filesystem, and recording said breaking of links in said persistent memory in a plurality of steps; and

altering said live filesystem to reflect changes associated with said breaking of
30 links, and recording said alteration in said persistent memory.

1/5

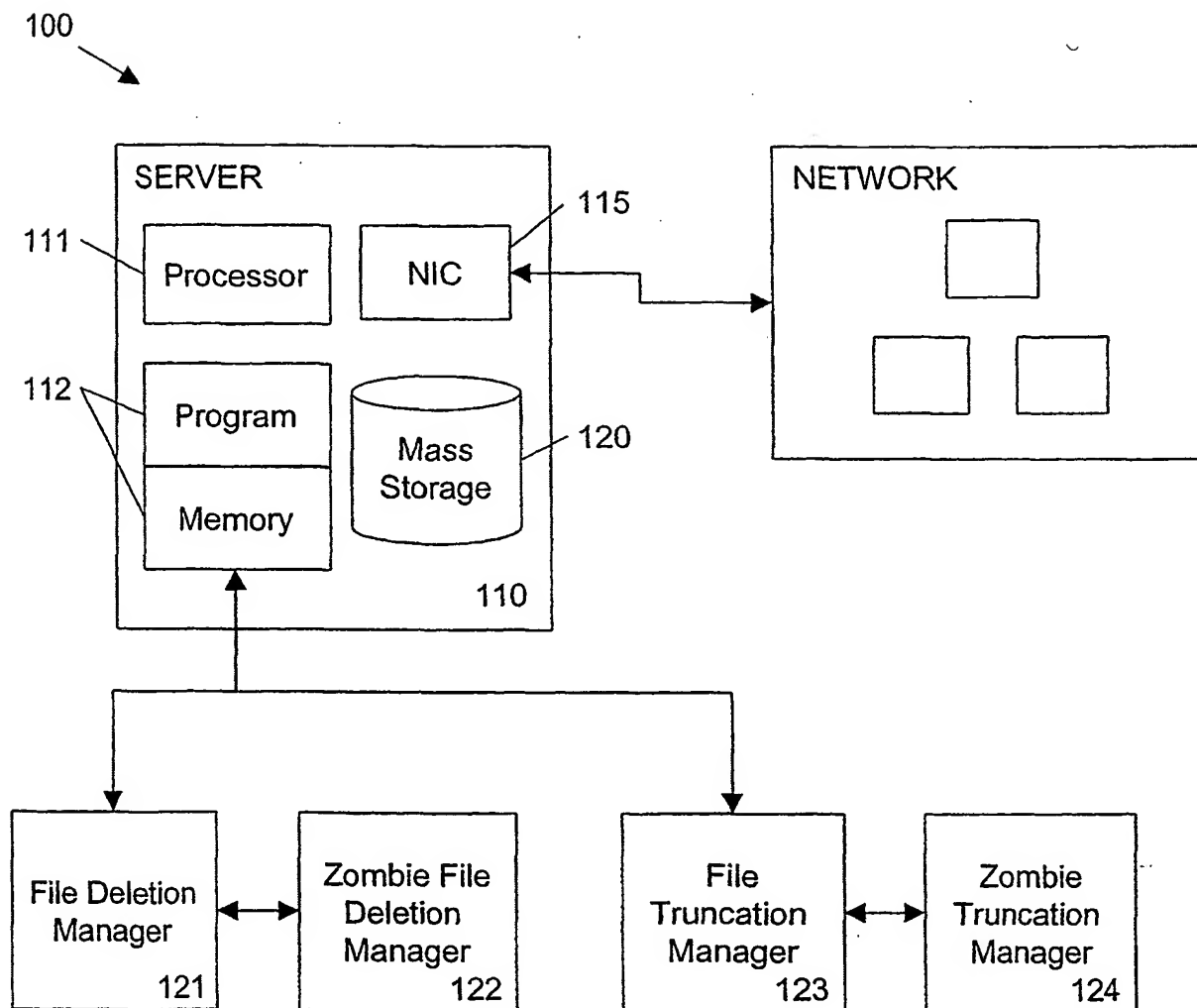


Fig. 1

2/5

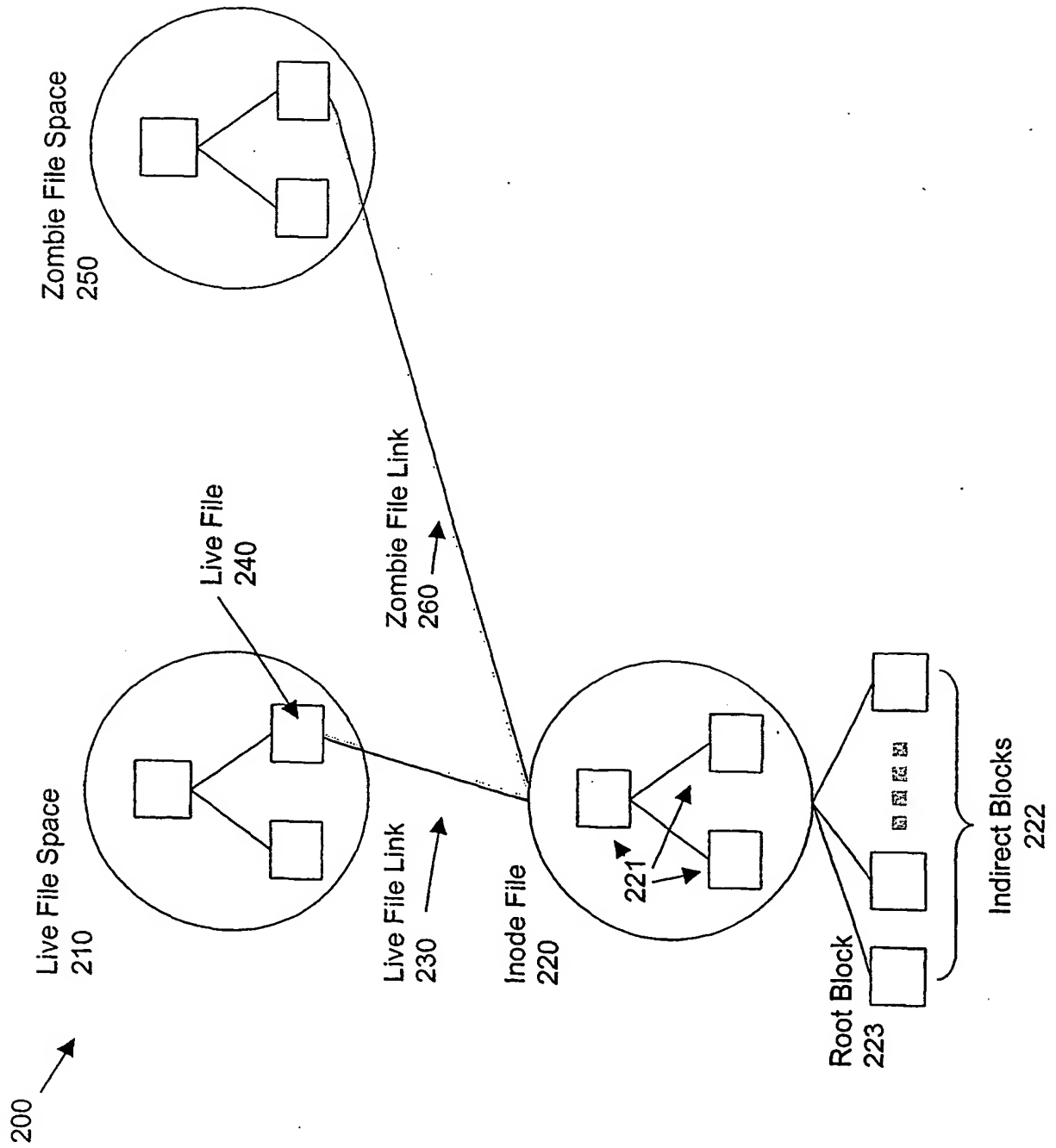


Fig. 2

3/5

Method 300

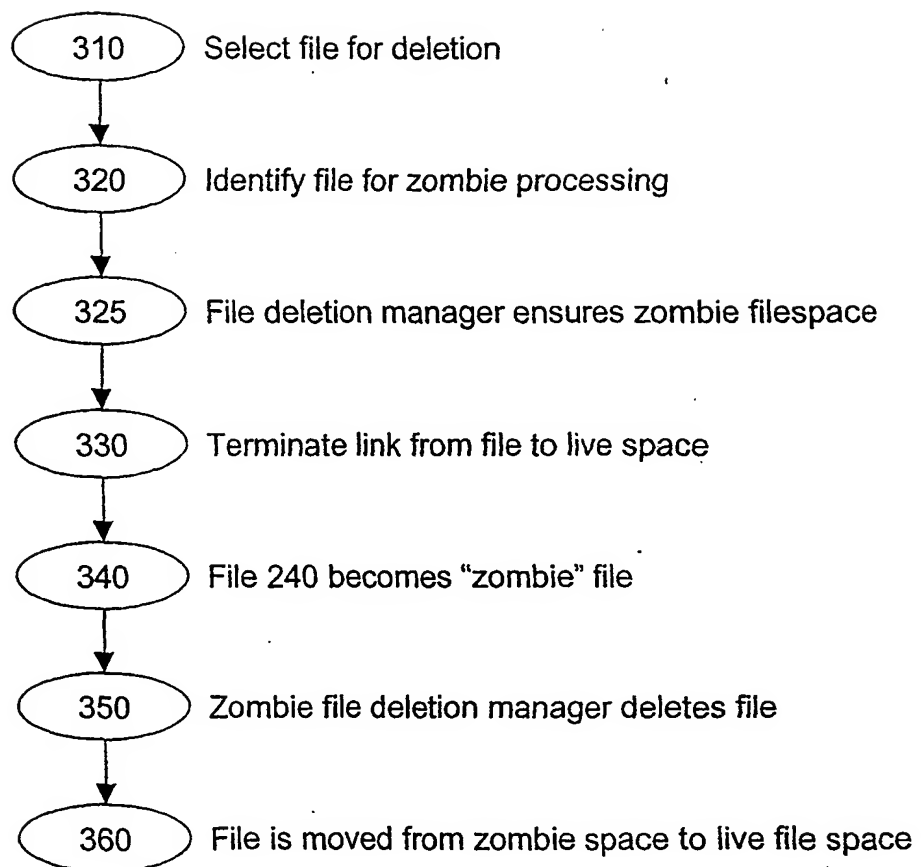


Fig. 3

4/5

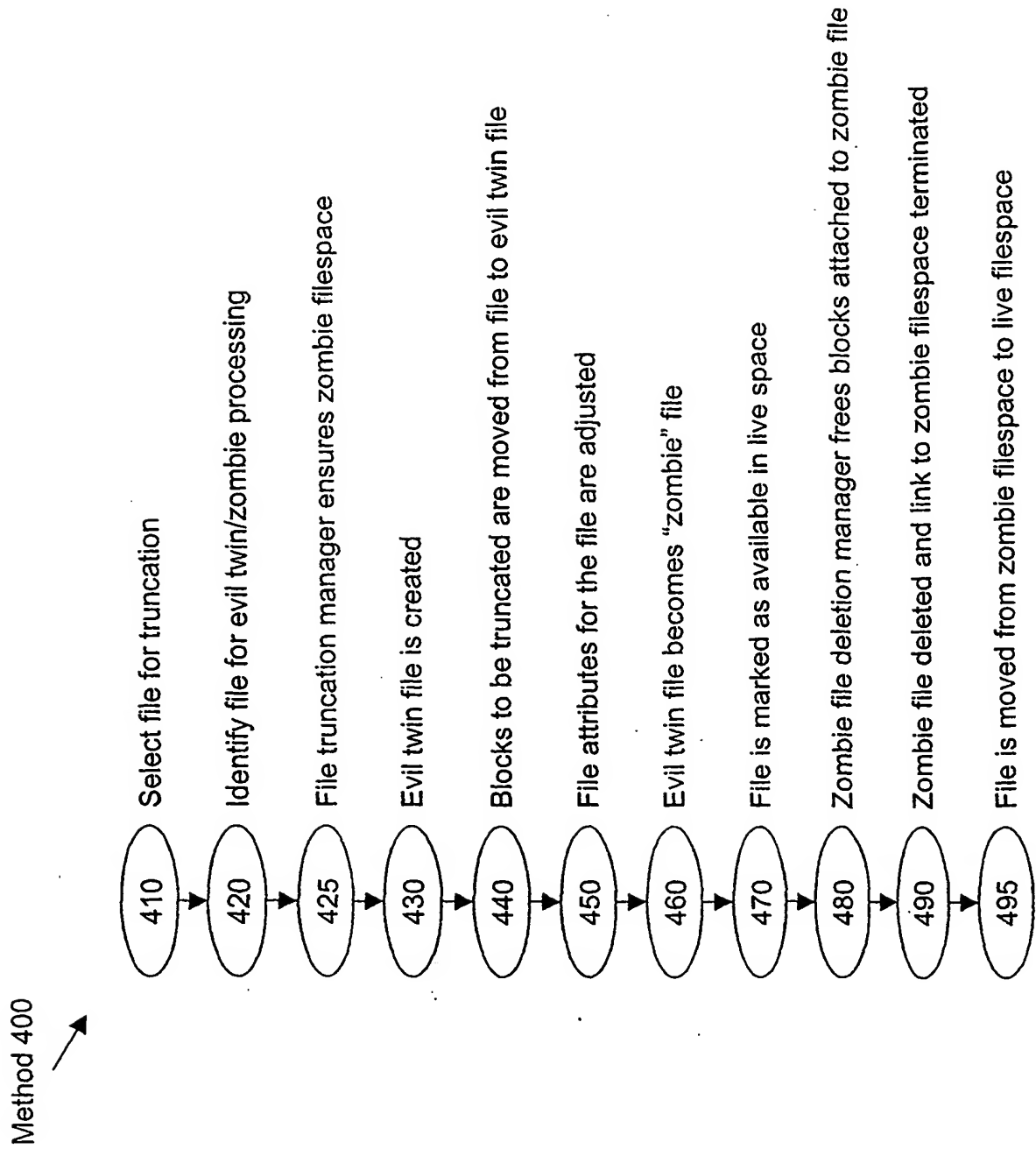


Fig. 4

5/5

Method 500

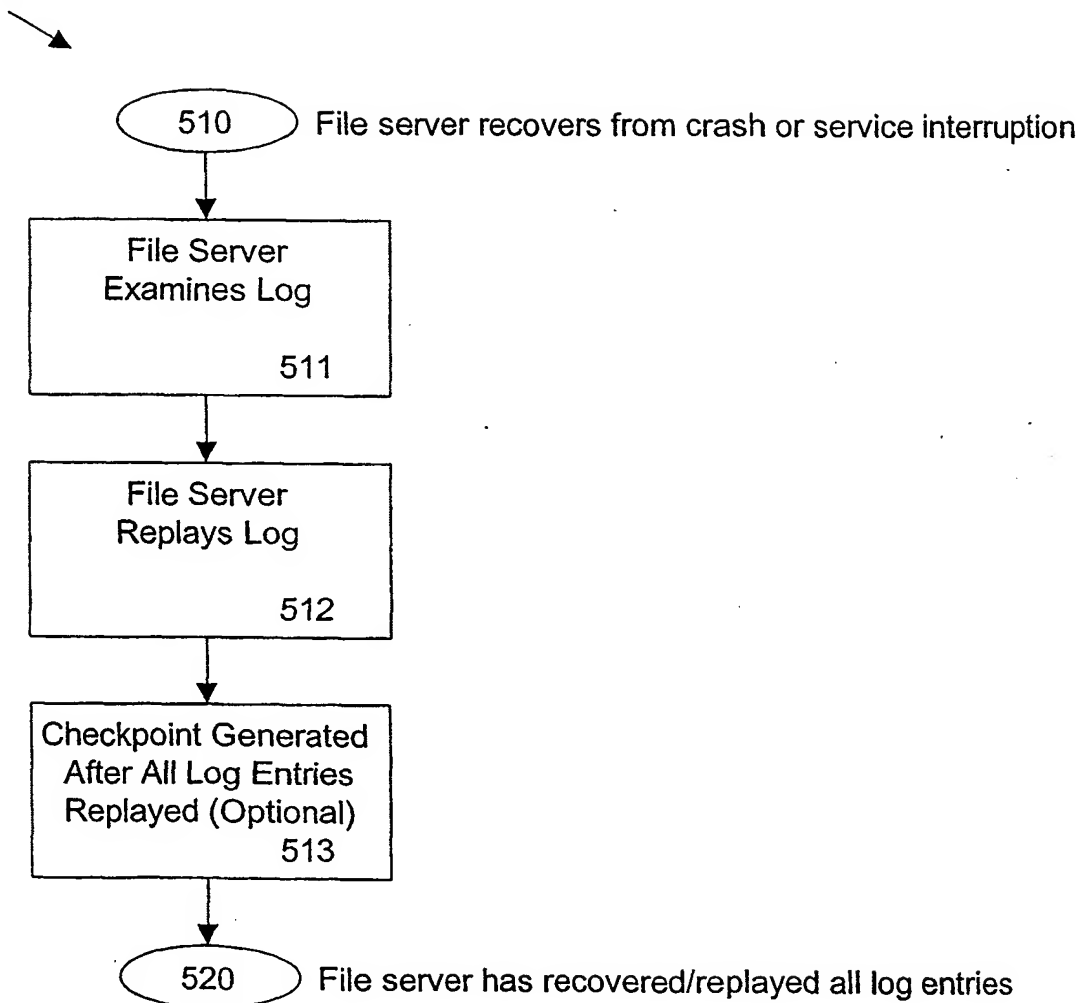


Fig. 5

THIS PAGE BLANK (USPTO)